

SOM에서 개체의 시각화

엄익현¹⁾, 허명희²⁾

요약

코호넨(T. Kohonen)의 자기조직화지도(Self-Organizing Map; SOM)은 저차원 그리드 공간에 고차원 다변량 자료를 축약하여 시각적으로 나타내는 비지도 학습법의 일종으로 최근 들어 통계 분석자들이 많은 관심을 가지고 있는 분야이다. 그러나 SOM은 개체공간의 연속형으로 표현되는 개체를 저차원 그리드공간에 승자노드에 비연속적으로 표현한다는 단점을 지니고 있다. 본 논문에서는 SOM을 통계적 목적으로 사용하기 위해 요구되는 그리드공간에 개체를 연속적으로 표현하는 방법들을 제안하고 활용 예를 제시하고자 한다.

주요용어: 자기조직화 지도(SOM), 코호넨(T. Kohonen), 시각화, IL-SOM, 부노드(k) SOM.

1. 연구 배경 및 목적

SOM은 1980년대 초반 핀란드의 전기공학자 코호넨(T. Kohonen)에 의해 개발된 비지도 학습(unsupervised learning) 신경망(neural network) 모형의 한 종류이다. 코호넨은 SOM의 특성을 시각화(visualization)과 축약화(abstraction)의 두 가지로 뽑았다 (Kohonen, 1998). 또한 SOM은 자기조직화(self organization)라는 과정을 통해 다차원 공간의 유사한 개체들을 서로 이웃하는 위치에 오도록 저차원 공간에 배치한다.

전통적인 SOM 알고리즘은 개체공간에 연속형으로 놓이는 개체를 저차원 그리드 공간에 비연속적으로 표현한다. 본 연구에서는 이 문제를 해결하기 위해 가능도(likelihood)를 이용하는 'IL-SOM'과 부노드를 이용하는 '부노드(k) SOM'을 제안하고 활용 예를 제시하고자 한다.

개체의 시각화에 대한 기존 연구로는 Goppert and Rosenstiel (1997), Campos and Carpenter (2000), 허명희 (2003)가 제안한 PC-SOM이 있다. 개체의 시각화에 대한 다른 방법으로 개체를 승자노드의 주변에 랜덤하게 분포 시키는 방법이 있다. 본 논문에서는 이 방법을 '임의 시각화 SOM'이라고 부르도록 한다. (엄익현, 2003).

2. 개체 시각화 방법의 제안

본 논문에서 제안하는 개체 시각화 방법은 개체벡터의 승자노드와 그 인접노드의 중량벡터에 대한 가능도(likelihood)를 이용하여 그리드 공간상에 개체벡터의 상대적인 위치를 표현하는 IL-SOM (Interpolating using Likelihood for SOM)과 승자노드 주변에 k^2 개의 부노드(subnode)를 배치하여 활용하는 부노드(k) SOM이다.

본 논문에서는 개체의 시각화에 대한 방법을 2차원으로 제한하여 제안하고자 한다. 그러나 1차원 또는 3차원 이상의 경우도 자연스럽게 구현 가능하다. 또한 그리드 구성이 사각형인 경우

1) (주) 지디에스코리아, 이사/통계학박사, (135-818) 서울시 강남구 논현동 81-10.

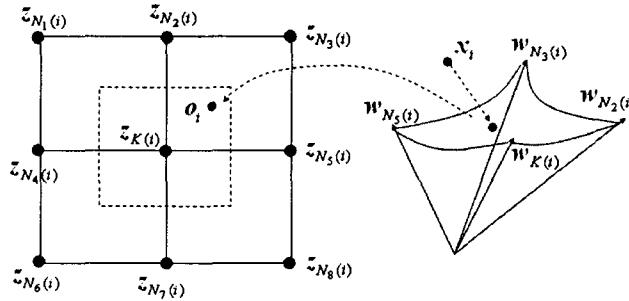
2) 고려대학교 통계학과, 교수, (136-701) 서울시 성북구 안암동 5가 1.

SOM에서 개체의 시각화

이든 육각형인 경우이든 동일한 방식으로 처리할 수 있으므로 사각형인 경우로 제한하여 다루고자 한다.

2.1 IL-SOM

SOM을 수행하면 p 차원 개체공간상의 m 개의 중량벡터 w_j 와 개체벡터 x_i 의 승자노드 $K(i)$ 의 중량벡터 $w_{K(i)}$ 를 얻게된다 ($i=1, \dots, n, j=1, \dots, m$). w_j 의 그리드 공간상의 이미지 좌표를 z_j 라고 하고, x_i 의 그리드 공간상의 이미지 좌표를 o_i 라고 하자. 그리고 $K(i)$ 를 둘러싸고 있는 8개의 인접노드를 $N_q(i)$ 로 표기하자 ($q = 1, \dots, 8$). 이를 그림으로 표현하면 그림 1과 같다.



(a) 2차원 그리드 공간 (b) p 차원 개체공간
그림 1: 저차원 그리드 공간에서의 개체 시각화

$L(w_j | x_i)$ 을 개체벡터 x_i 의 중량벡터 w_j 에 대한 가능도로 다음과 같이 정의하자.

$$L(w_j | x_i) = \text{상수} \cdot \exp[-\frac{1}{2}(x_i - w_j)^T \Sigma^{-1} (x_i - w_j)] \quad (1)$$

여기서, $j = \{K(i), N_1(i), \dots, N_8(i)\}$ 이고, Σ 는 알려지지 않은 상수 β 에 대해서 $\Sigma = \beta I$ 라 가정하자. 여기서, 가능도에 대한 Gauss 분포를 가정하였지만 특별한 이유가 있는 것은 아니다. 그러나 Gauss 분포로써 평균과 산포를 동시에 고려할 수 있기 때문에 쉽게 생각할 수 있는 한 방법이다.

o_i 를 구하기에 앞서 승자노드를 중심으로 인접노드에 의해서 구성되는 네 개의 사각형 중 개체가 표현될 방향을 결정한다. 방향을 결정하는 방법은 각 방향을 둘러싸고 있는 네 개의 중량벡터의 가능도들의 합을 구한 후 가장 큰 가능도합을 가지는 방향을 개체가 표현될 방향으로 결정한다.

이제 o_i 를 구해 보자. 그림 1에서 결정된 방향이 승자노드의 우측상단 방향인 경우를 예를 들어 설명하면, x_i 의 승자노드와 인접노드의 중량벡터에 대한 네 개 가능도들의 합을 구하고, 이로부터 각각의 중량벡터가 차지하는 비율 p_j 를 다음과 같이 구한다. 물론 $\sum_j p_j = 1$ 이다.

$$p_j = \frac{\exp[-1/(2\beta)(x_i - w_j)^T (x_i - w_j)]}{\sum_{\{k: K(i), N_2(i), N_3(i), N_5(i)\}} \exp[-1/(2\beta)(x_i - w_k)^T (x_i - w_k)]} \quad (2)$$

네 개의 p_j 값을 이용하여 다음과 같이 o_i 를 구함으로써 개체를 그리드 공간상에 표현한다. 즉,

$$o_i = \sum_{(j: K(i), N_2(i), N_3(i), N_5(i))} p_j \cdot z_j. \quad (3)$$

2.2 가장자리 노드가 승자노드인 경우의 개체의 시각화

그리드공간상에 개체를 표현할 때 그리드 바깥쪽에 가상의 노드를 설정하여 가장자리 노드의 바깥쪽에도 개체가 표현될 수 있게 하는 것이 바람직하다.

가상의 노드의 개체공간상의 중량벡터를 w_{e1} 라고 하고 가장자리 노드의 중량벡터를 w_1 , w_1 을 중심으로 w_{e1} 의 정반대 방향 인접노드의 중량벡터를 w_2 라 하자. 선형보간법에 의하여 w_{e1} 를 설정하면 다음과 같다.

$$w_{e1} = 2w_1 - w_2 \quad (4)$$

마찬가지 방식으로, 바깥쪽 가상 노드 중 4개의 모퉁이에 위치한 가상 노드는 다음과 같이 설정한다.

$$w_{e2} = 4w_1 - 2w_2 - 2w_3 + w_4 \quad (5)$$

그림 2는 중량벡터 w_{e1} 과 w_{e2} 에 대응하는 그리드 공간상의 노드 z_{e1} 과 z_{e2} 를 보여준다.

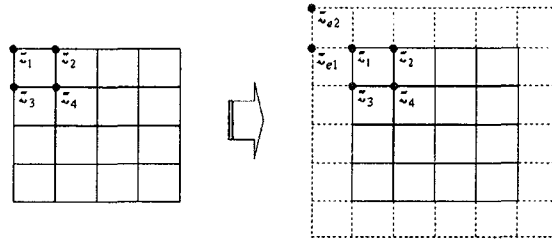


그림 2: 그리드의 바깥쪽에 가상의 노드 설정

2.3 상수 β 의 결정

먼저 o_i 를 개체공간 상의 네 개의 중량벡터에 의해서 구성되는 영역에 옮겨 놓는 문제를 생각해 보자. 그림 3(a)와 같은 상황을 설정하자. 2차원 그리드 공간에서 o_i 은 승자노드 $z_{K(i)}$ 와 가로방향의 인접노드 $z_{N_5(i)}$ 를 $a:(1-a)$ 로 내분하는 위치에 있다. 또한 승자노드 $z_{K(i)}$ 와 세로방향의 인접노드 $z_{N_2(i)}$ 를 $b:(1-b)$ 로 내분하는 위치에 있다. 따라서 다음과 같은 관계가 성립한다.

$$o_i = (1-a)(1-b)z_{K(i)} + (1-a)b z_{N_2(i)} + a(1-b)z_{N_5(i)} + ab z_{N_3(i)} \quad (6)$$

식 (6)의 관계를 개체공간으로 옮겨가면 다음과 같이 표현할 수 있는데, 이를 통해서 개체공간상의 네 개의 중량벡터로 구성되는 영역에 o_i 의 이미지 x_i^w 를 표현할 수 있다. 그림 3(b)을 보라.

$$x_i^w = (1-a)(1-b)w_{K(i)} + (1-a)b w_{N_2(i)} + a(1-b)w_{N_5(i)} + ab w_{N_3(i)} \quad (7)$$

SOM에서 개체의 시각화

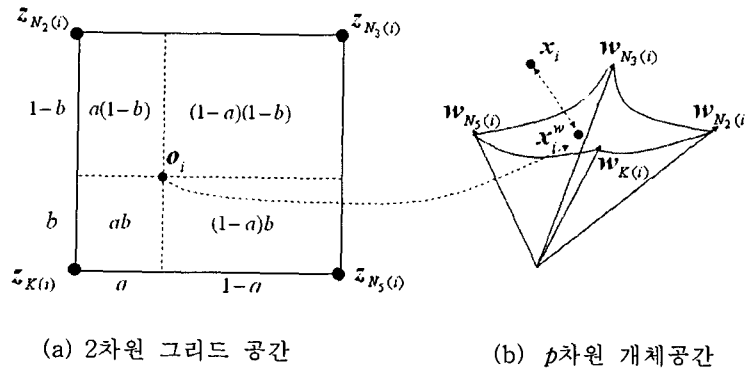


그림 3: 개체공간에 o_i 의 이미지 x_i^w 의 표현

다음으로 개체벡터 x_i 와 x_i^w 간의 제곱거리의 합을 다음과 같이 정의하고, 제곱거리의 합을 가장 작게 하는 β 을 선택하기로 한다. x_i^w 가 β 의 함수라는 의미에서 $x_i^w(\beta)$ 로 표현했다.

$$Q_\beta = \min_{\beta} \sum_{i=1}^n \|x_i - x_i^w(\beta)\|^2 \quad (8)$$

이 때 구한 최소제곱거리를 Q_β 로 표기하고 이를 $r \times c$ 그리드에서의 개체표현지수라 명명한다.

2.4 변수의 표현

SOM 학습을 통해서 얻어진 저차원 그리드 공간에 변수들을 표현하는 방법을 살펴보자. j 번째 변수의 방향을 크기 $p \times 1$ 벡터 v_j 로 나타내자 ($j = 1, \dots, p$). SOM은 비선형사영이므로 변수 축이 저차원공간에 하나의 직선이 아닌 곡선형태로 표현된다. 따라서 -3 표준편차에서 +3 표준편차까지의 범위를 0.5 표준편차 간격으로 다음과 같이 13개의 벡터를 표현하는 것을 제안한다.

$$v_j = (0, \dots, s, \dots, 0) \quad (9)$$

여기서, $s = -3, -2.5, \dots, 0, \dots, +2.5, +3$.

식 (9)의 변수벡터 v_j 를 그리드 공간에 표현하는 방법은 다음과 같이 요약된다.

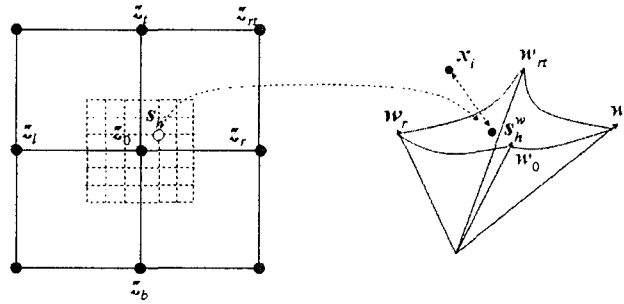
- [0] $s = -3$ 으로 설정한다.
- [1] v_j 의 승자노드를 구한다.
- [2] v_j 의 그리드 공간상의 좌표 o_{v_j} 가 놓일 방향을 선택한다.
- [3] 식 (2)와 (3)에서와 같은 방법으로 o_{v_j} 를 구한다.
- [4] s 를 0.5 단위로 증가하여 3이 될 때까지 [1]~[3]을 반복한다.

2.5 부노드(k) SOM

그리드 공간에서 승자노드를 z_0, z_0 를 둘러싸고 있는 네 개의 인접 노드를 z_l, z_r, z_t, z_b 라고 하고 이들 노드에 대한 개체공간의 증량벡터를 각각 w_0, w_b, w_r, w_t, w_b 라고 하자. 그림 4를 보라.

그리드 공간의 노드 z_0 와 z_l, z_0 와 z_r, z_0 와 z_t, z_0 와 z_b 를 연결하는 네 개의 선분을 k

등분하는 선분을 각각의 선분에 수직으로 긋는다 (여기서 k 는 홀수). 이들 선분에 의해서 만나는 점을 모두 $(2k-1)^2$ 개 얻게 된다. 이들 점 중에서 주변노드 z_l, z_r, z_t, z_b 보다 승자노드 z_0 에 더 가까운 점을 k^2 개 얻을 수 있는데 이를 부노드(subnode)라 부르고 s_h 라 표기하자 ($h = 1, \dots, k^2$).



(a) 2차원 그리드 공간 (b) p 차원 개체공간

그림 4: 부노드(k) SOM을 이용한 개체의 표현

부노드(k) SOM은 이들 k^2 개의 부노드들을 2.3절에서 제안한 IL-SOM의 식 (6)과 같은 방법으로 개체공간에 표현한다. 개체공간에 표현된 부노드를 s_h^w 로 표기하자.

부노드 s_h 를 개체공간에 표현한 후 개체벡터 x_i 와 s_h^w 와의 제곱거리를 구한다. 이를 k^2 개 부노드에 대해서 모두 실시한 후 최소제곱거리를 갖는 s_h^w 에 대응하는 부노드 s_h 를 개체 x_i 의 그리드 공간상의 좌표로 정하는 것을 제안한다. 즉,

$$\min_h \|x_i - s_h^w\|^2, \quad h = 1, \dots, k^2. \quad (10)$$

3. 붓꽃 자료 예

본 논문에서 제안한 IL-SOM과 부노드(k) SOM을 코호넨의 SOM, Goppert and Rosenstiel (1997)가 제안한 ISOM과 임의의 시각화 SOM 등과 함께 Fisher의 붓꽃 자료에 적용해본 후 다섯 가지 방법의 개체표현의 수행능력(performance)을 비교하고자 한다.

Fisher의 붓꽃 자료(Iris Data)는 세 가지 품종(1: setosa, 2:versicolor, 3:verginica)의 붓꽃으로부터 각각 50개, 총 150개의 개체를 추출하여 측정된 네 개의 변수 X_1 : 꽃받침 조각의 길이 (sepal length), X_2 : 꽃받침 조각의 폭 (sepal width), X_3 : 꽃잎의 길이 (petal length), X_4 : 꽃잎의 폭 (petal width)으로 구성되어 있다. SOM을 수행하기 위한 입력변수는 $X_1 \sim X_4$ 이고 사전에 표준화 변환되었다 (평균 0, 표준편차 1).

SOM 수행시, 사전에 그리드 크기를 5x5, 초기 학습률 0.25, 최종 학습률 0.001, 초기 주변거리 2, 최종 주변거리 1, 주기 당 반복 50회 등으로 지정한 결과이다. 연구자가 작성한 SAS/IML 프로그램을 SOM 산출을 위하여 사용하였다. 각 그림의 노드 위의 큰 숫자는 해당노드를 승자로 가지는 개체의 수이며, 그리드상에 넓게 퍼져있는 작은 숫자는 품종을 표시하며 해당 위치가 개체의 표현위치이다.

SOM에서 개체의 시각화

그림 5는 IL-SOM ($\beta=0.09$), 부노드(k) SOM, ISOM과 임의시각화 SOM에 대한 결과를 보여준다. 여기서 ISOM의 경우 $\lambda = 0.001$ 로 놓았다.

Fisher의 붓꽃 자료에 대해서 본 논문에서 제안한 두 가지 방법과 코호넨의 SOM, ISOM, 그리고 임의 시각화 SOM의 개체표현지수, 즉 제곱거리의 합을 비교한 결과를 표 1에 정리하였다.

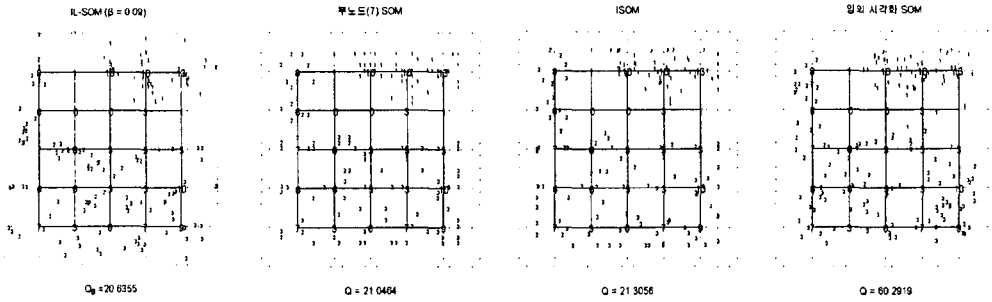


그림 5: 붓꽃 자료에 대한 5x5 그리드 공간에서 개체표현

표 1: 붓꽃 자료에 대한 5x5 그리드 공간에서 각 방법의 개체표현 결과 비교

방법	최적의 β	Q_β (또는 Q)	비교
IL-SOM	0.09	20.6355	가장 작음
부노드(7) SOM	-	21.0464	
코호넨 SOM	-	38.7422	
ISOM ($\lambda = 0.001$)	-	21.3056	
임의 시각화 SOM	-	60.2919	

4. 모의자료 예

이 절에서는 구조가 잘 알려진 모의자료에 대해서 그리드공간상에 개체와 변수를 표현해 보도록 하자. 개체 및 변수의 표현이 기대하는 바와 일치하는 가를 확인하는 것이 목적이다.

모의자료는 네 개의 그룹으로 구분되며, 각 그룹별 50개의 개체, 총 200개의 개체를 생성하였다. 각 그룹은 다음과 같은 평균과 공분산행렬을 가지는 다변량정규분포를 따른다.

$$1 \text{ 그룹 : 평균이 } \left(-2.5, -\frac{2.5 \cdot \sqrt{3}}{3}, -\frac{2.5 \cdot 2 \cdot \sqrt{2}}{4 \cdot \sqrt{3}}, 0 \right) \text{이고 공분산행렬이 } I_4,$$

$$2 \text{ 그룹 : 평균이 } \left(2.5, -\frac{2.5 \cdot \sqrt{3}}{3}, -\frac{2.5 \cdot 2 \cdot \sqrt{2}}{4 \cdot \sqrt{3}}, 0 \right) \text{이고 공분산행렬이 } I_4,$$

$$3 \text{ 그룹 : 평균이 } \left(0, \frac{2.5 \cdot 2 \cdot \sqrt{3}}{3}, -\frac{2.5 \cdot 2 \cdot \sqrt{2}}{4 \cdot \sqrt{3}}, 0 \right) \text{이고 공분산행렬이 } I_4,$$

$$4 \text{ 그룹 : 평균이 } \left(0, 0, \frac{2.5 \cdot 6 \cdot \sqrt{2}}{4 \cdot \sqrt{3}}, 0 \right) \text{이고 공분산행렬이 } I_4.$$

이 자료의 네 그룹의 평균을 연결하면 한 변의 길이(유클리드거리)가 5이고, 무게중심이 원점인 정사면체를 얻게 된다. 여기서 변수 X_4 는 그룹의 구분과는 상관이 없게 생성된 의미없는 변수이다. 이러한 정사면체 구조를 가지는 자료의 경우 2차원 평면에 선형사영을 하였을 때 어느 평면에 사영을 하던 4개의 그룹이 잘 구분이 되지 않고 일부 그룹은 겹쳐서 표현된다. 이러한 문제가 발생하는 이유는 정사면체가 가지는 위상적인 구조가 선형사영에는 알맞지 않다는

데 기인한다. 따라서 이 경우에는 비선형사영을 고려하는 것이 더 타당해 보이며 SOM을 통해 비선형사영을 시도하고자 한다.

그림 6은 모의 자료에 대한 그리드공간에 변수를 표현한 것으로 (a)~(d)는 각각 $X_1 \sim X_4$ 의 변수를 표현한 그림이다. IL-SOM ($\beta = 0.43$)을 이용하여 5×5 그리드공간에 변수를 표현하였다. 변수의 표현이 기대했던 바와 일치하지를 알아보기 위하여 그림 6에서 각 변수를 살펴보자. 변수 X_4 은 군집화에 별다른 영향을 미치지 못하고 있으며, 변수 X_1 은 1 군집과 2 군집의 대조, X_2 는 1,2 군집과 3 군집의 대조, X_3 는 1,2 군집과 4 군집의 대조를 의미는 것으로 나타나고 있다. 즉, 그림 6에서 기대했던 바와 같이 모의자료의 변수를 표현해 주고 있다.

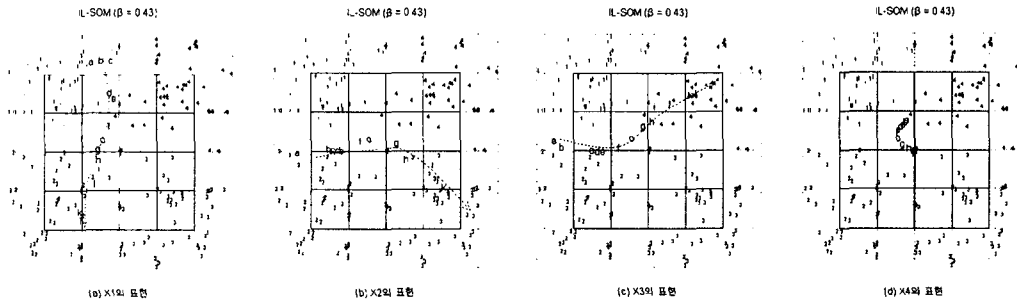


그림 6: 모의자료에 대한 5×5 그리드공간에서 IL-SOM을 이용한 변수표현
(a: -3, b: -2.5, c: -2, d: -1.5, e: -1, f: -0.5, o: 0, g: 0.5, h: 1, i: 1.5, j: 2, k: 2.5, l: 3)

5. 맺음 말

본 연구에서는 SOM에서 개체를 그리드공간상에 시각적으로 표현할 수 있는 두 가지 방법을 제안하였다. 이를 통해 전통적인 SOM이 가지는 이산형 출력에 대한 모순을 개선하고, 개체를 저차원 그리드공간에 표현하는 것이 가능하여 결과의 시각적 해석 및 변수의 표현을 통해 각 군집의 특성 파악이 가능하게 하였다.

참고문헌

[1] 엄익현 (2003). 코호넨 자기조직화지도(SOM)의 통계적 활용. 고려대학교 대학원 통계학과 박사학위 논문.
 [2] 허명희 (2003). "주성분 자기조직화 지도 PC-SOM", 응용통계연구 16권 2호. 321-334.
 [3] Campos, M.M., and Carpenter, G.A. (2000). "Building adaptive basis functions with a continuous self-organizing map," Neural Processing Letter, 11. 59-78.
 [4] Goppert, J. and Rosenstiel, W. (1997). "The continuous interpolating self-organizing map," Neural Processing Letter, 5. 185-192.
 [5] Kohonen, T. (1995). Self-Organizing Maps. Springer, Berlin.
 [6] Kohonen, T. (1998). "The self-organizing map," Neurocomputing 21, 1-6.